

(10) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2000-20501

(P2000-20501A)

(43) 公開日 平成12年1月21日(2000.1.21)

(5) Int.Cl. <sup>1</sup>	識別記号	P I	フット* (参考)
G 0 6 F 17/10		G 0 6 F 15/31	2 5 B 0 4 5
15/18	3 9 9	15/18	3 9 0 Z 5 B 0 6 6

審査請求 未請求 請求項の数7 O L (全 17 頁)

(21) 出願番号 特願平10-188849  
(22) 出願日 平成10年7月3日(1998.7.3)

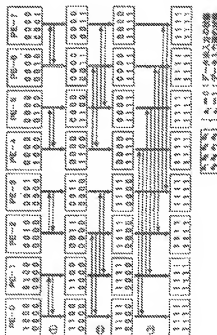
(71) 出願人 000003078  
株式会社東芝  
神奈川県川崎市幸区瀬川町72番地  
(72) 発明者 上松 幹夫  
神奈川県横浜市磯子区新杉田町8番地 株式会社東芝横浜事業所内  
(74) 代理人 100083161  
弁護士 外川 亮明  
Fターム(参考) 58045 A67 B302 B628 B647 G002  
G012  
58050 A04A B045 B034 F006

(54) 【発明の名称】 並列計算機システム及びその演算処理装置間の通信方法

(57) 【要約】

【課題】 並列計算機の各演算処理装置で分散処理されたデータを効率的に集結する。

【解決手段】 識別番号0, 1, ..., 2<sup>n</sup>-1が付与された2<sup>n</sup>台の演算処理装置と識別記号装置及び通信手段を備えた並列計算機システムで、2<sup>n</sup>個の小配列に分割して各演算処理装置に分配／演算処理されたデータ配列を1つの配列に集結する際に、識別番号Nに對して2進法で表した識別番号Nの2<sup>i</sup>の位の値を反転させた番号N'を算出させ、識別番号Nの演算処理装置と識別番号N'の演算処理装置の間でデータ配列の演算処理結果を相互に送受信する操作；をi=0からi=n-1まで順次行う。この際、jとなるjに對しては、操作jの際に、識別番号N, N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作(j-1)までで得られた演算処理結果も送受信する。



1

【特許請求の範囲】

【請求項1】 固有の識別子を有する少なくとも2台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、2台の小配列に分割して2台の演算処理装置に分配され各演算処理装置で演算処理されたデータ配列を再び1つの配列に集結する際に、2台の演算処理装置に識別番号0、1、…、 $2^k - 1$ を付与し、識別番号Nの演算処理装置に対して2進法で表した識別番号Nの2<sup>k</sup>の位の数を反転させた番号N'を識別番号とす演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受する操作1をi = 0からi = n - 1まで順次行い、i > nになるに対しては、操作1の際に、識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加え操作(j - 1)までで得られた演算処理結果を送受することにより2台の演算処理装置間でj回の操作でデータ配列を集結させることを特徴とする並列計算機システム。

【請求項2】 固有の識別子を有する(2<sup>k</sup> + k)台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、(2<sup>k</sup> + k)箇の小配列に分割して(2<sup>k</sup> + k)台の演算処理装置に分配、演算処理されたデータ配列を再び1つの配列に集結する際に、前記(2<sup>k</sup> + k)台の演算処理装置に個別記憶手段及び通信手段を備えた(2<sup>k</sup> - k)台の演算処理装置を加えた2<sup>k+1</sup>台からなる演算処理装置群を形成し、この演算処理装置群を構成する2<sup>k+1</sup>台の演算処理装置に識別番号0、1、…、2<sup>k+1</sup> - 1を付与し、識別番号Nの演算処理装置に対して2進法で表した識別番号Nの2<sup>k</sup>の位の数を反転させた番号N'を識別番号とす演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受する操作1をi = 0からi = mまで順次行い、i > mになるに対しては、操作1の際に、N ≤ 2<sup>k</sup> + kなる識別番号Nの演算処理装置からj回目の演算処理結果の演算処理結果を送受することにより(2<sup>k</sup> + k)台の演算処理装置において(m + 1)回の操作でデータ配列を集結させることを特徴とする並列計算機システム。

【請求項3】 固有の識別子を有する(2<sup>k</sup> + k)台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機

2

システムにおいて、(2<sup>k</sup> + k)箇の小配列に分割して(2<sup>k</sup> + k)台の演算処理装置に分配、演算処理されたデータ配列を再び1つの配列に集結する際に、この(2<sup>k</sup> + k)箇のデータ配列に(2<sup>k</sup> - k)箇の空の小配列を追加することで前記データ配列を小配列2<sup>k+1</sup>箇分の配列に拡張し、前記(2<sup>k</sup> + k)台の演算処理装置に、個別記憶手段及び通信手段を備えた(2<sup>k</sup> - k)台の演算処理装置を加えた2<sup>k+1</sup>台からなる演算処理装置群を形成し、この演算処理装置群を構成する2<sup>k+1</sup>台の演算処理装置に識別番号0、1、…、2<sup>k+1</sup> - 1を付与し、識別番号Nの演算処理装置に対して2進法で表した識別番号Nの2<sup>k</sup>の位の数を反転させた番号N'を識別番号とす演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受する操作1をi = 0からi = mまで順次行い、i > mになるに対しては、操作1の際に、識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作(j - 1)までで得られた演算処理結果を送受することにより(2<sup>k</sup> + k)台の演算処理装置において(m + 1)回の操作でデータ配列を集結させることを特徴とする並列計算機システム。

【請求項4】 n > mなるn、mについて、固有の識別子を有する(2<sup>k</sup> + 2<sup>k</sup>)台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、(2<sup>k</sup> + 2<sup>k</sup>)箇の小配列に分割して(2<sup>k</sup> + 2<sup>k</sup>)台の演算処理装置に分配、演算処理されたデータ配列を再び1つの配列に集結する際に、前記(2<sup>k</sup> + 2<sup>k</sup>)台の演算処理装置を2<sup>k</sup>台からなるグループA、と2<sup>k</sup>台からなるグループBに分割し、また前記データ配列を初めの2<sup>k</sup>箇の小配列からなる配列A、とその後の2<sup>k</sup>箇の小配列からなる配列A'、の2つに分割し、この配列A、A'、をそれぞれグループA、とA'、と対応づけて分配、演算処理を行い、グループA、の2<sup>k</sup>台の演算処理装置に識別番号0、1、…、2<sup>k</sup> - 1を付与し、識別番号Nの演算処理装置に対して2進法で表した識別番号Nの2<sup>k</sup>の位の数を反転させた番号N'を識別番号とす演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受する操作1をi = 0からi = n - 1まで順次行い、i > nになるに対しては、操作1の際に識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作(j - 1)までで得られた演算処理結果を送受することによりグループA、内で配列A、を連結させる第1の工程と、グループB、の2<sup>k</sup>台の演算処理装置に識別番号0、1、…、2<sup>k</sup> - 1を付与し、識別番号Nの演算処理装置に対して2進法で表した識別番号Nの2<sup>k</sup>の位の数を反転させた番号N'を鑑

別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受信する操作（i）をi=0からi=q-1まで順次行い、j>0なるjに対して、操作jの際に識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作（j-1）までで得られた演算処理結果を送受信することによりグループG<sub>1</sub>、内記配列A、を連結させる第2の工程と、グループG<sub>1</sub>からグループG<sub>2</sub>の各演算処理装置に配列A<sub>1</sub>を、グループG<sub>2</sub>からグループG<sub>3</sub>の各演算処理装置に配列A<sub>2</sub>を送信する第3の工程とを有し、第1の工程と第2の工程を並列に実行した後に第3の工程を行なうことにより（2' + 2''）台の演算処理装置においてデータ配列を連結させることを特徴とする並列計算機システム。

【請求項5】 固有の識別子を有する複数の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備えた並列計算機システムにおいて、

【数1】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

個のデータ配列（但し、 $n_1 > n_2 > n_3 > \dots > n_k \geq 0$ ）に分割して

【数2】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

台の演算処理装置に分配・演算処理されたデータ配列を再び1つの配列に連結する際に、これらの演算処理装置のうち

【数3】

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

台をそれぞれグループG<sub>1</sub>、G<sub>2</sub>、…、G<sub>k</sub>としてk個のグループに分割するとともに、前記データ配列のうち

【数4】

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

個の小配列をそれぞれ配列A<sub>1</sub>、A<sub>2</sub>、…、A<sub>k</sub>としてk個の配列に分割し、このk個の配列とk個のグループG<sub>1</sub>、G<sub>2</sub>、…、G<sub>k</sub>とを1対1に対応づけて分配・演算処理を行い、1台よりなる各小グループG<sub>1</sub>、G<sub>2</sub>、…、G<sub>k</sub>の（2のj乗）台の演算処理装置に識別番号0、1、…を付し、識別番号Nの演算処理装置に対し2進法で変換した識別番号Nの2<sup>j</sup>の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受信する操作（i）をi=0からi=q-1まで順次行い、j>0なるjに対して、操作jの際に識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作（j-1）までで得られた演算処理結果を送受信する

ことによりグループG<sub>1</sub>、内の演算処理装置でデータ配列A<sub>1</sub>を連結させるグループ内工程pを実行し、グループ内工程（k-1）が終了した後、グループG<sub>1</sub>の演算処理装置から配列A<sub>1</sub>の演算結果をグループG<sub>2</sub>の演算処理装置に送信するグループ間工程kを実行し、次に、グループG<sub>2</sub>の各演算処理装置に供給された配列A<sub>2</sub>の演算結果を、グループG<sub>2</sub>の演算処理装置からq>0となる全てのqに対しグループG<sub>2</sub>に属する各演算処理装置に送信するとともに、グループG<sub>2</sub>の演算処理装置から、グループG<sub>2</sub>自身の演算結果である配列A<sub>2</sub>、及びグループG<sub>2</sub>の演算処理装置から受信した配列A<sub>1</sub>、…、A<sub>1</sub>の演算結果をグループG<sub>3</sub>の演算処理装置に送信するグループ間工程p、p=k-1からp=2までpに閉じて降順に実行することにより、

【数5】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

台の演算処理装置においてデータ配列を連結させることを特徴とする並列計算機システム。

20 【請求項6】 請求項5または前記記載の並列計算機システムを用いて演算処理装置のグループ間でデータ交換を行う場合、 $p > q$ なるp、qについて、2'台の演算処理装置からなるグループG<sub>1</sub>で連結され共有されているデータ配列Aと、2''台の演算処理装置からなるグループG<sub>2</sub>で連結され共有されているデータ配列Bとを、グループG<sub>1</sub>、G<sub>2</sub>間で相互に送受信する際に、グループG<sub>1</sub>のなかから選択される2'台の演算処理装置をグループG<sub>2</sub>の各演算処理装置と1対1に対応させてグループG<sub>2</sub>の各演算処理装置にデータ配列Aを送信する操作を並列に実施するとともに、グループG<sub>2</sub>を、それぞれが2''台の演算処理装置からなる小グループα<sub>1</sub>、α<sub>2</sub>、…、α<sub>r</sub>（r=2'）に分割して、各々の小グループとグループG<sub>1</sub>のi'台の各演算処理装置とを1対1に対応させ、小グループα<sub>i</sub>のなかから選択されるi'台の演算処理装置に対して、小グループα<sub>i</sub>に対応するグループG<sub>1</sub>の演算処理装置からデータ配列Bを送信した後、小グループα<sub>i</sub>の演算処理装置間でデータ配列Bを送受信する操作（i）をi=1からi=rに閉じて並列に実行するすることにより、2'台の演算処理装置と2''台の演算処理装置にデータ配列Aとデータ配列Bを共有させることを特徴とする並列計算機システム。

【請求項7】 請求項1乃至6のいずれか記載の並列計算機システムを用いて2台の演算処理装置間でデータを交換する工程は、演算処理装置の識別番号の大きい方から小さい方にデータを送る第1の送信工程と、演算処理装置の識別番号の小さい方から大きい方にデータを送る第2の送信工程とからなり、この第1の送信工程と第2の送信工程のうちから選択される1工程を先行した後に、続いて他の1工程を行うことを特徴とする並列演算機システムの演算処理装置間の通信方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、通信手段および制御装置を備えた多数の演算処理装置からなり、特に並列計算を目的とした並列計算システム及びその演算処理装置の通信方法に関する。

【0002】

【従来の技術】原子力施設をはじめとする大規模な施設の設計においては、例えば遮蔽設計などにおける放射線挙動計算、炉心設計における炉心性状平衡解析などの大規模な計算がかなりの頻度で要求される。この要求に応えるためには大規模な計算速度の向上が必要である。このため最近では、通信手段と制御の配役装置を備えた多数の演算処理装置を用いて、1台の演算処理装置しか持たない計算機を使用していたのでは得られないような高速度で、解析を行うことが考察されている。

【0003】例えば炉心設計であれば、原子炉の炉心を複数の燃料集合体からなる幾つかのセグメントに分割し、それぞれのセグメントを1つの演算処理装置に対応させて、出力計算と熱水力計算を各々の演算処理装置で並列に計算させる。セグメント間での中性子束の流出入およびチャンネル間の冷却材の圧力バランスを解析する際には、前記通信手段によりセグメント境界の中性子

\*束、各チャンネルの圧力損失のデータを演算処理装置間でやり取りすることで、空間的に連続した解析が行われる。

【0004】また、遮蔽設計であれば、例えば原子炉の炉心、冷却材、遮蔽体などを含む全体系を幾つかの小領域に分割し、それぞれの小領域を1つの演算処理装置に対応させて、放射線束分布計算を各々の演算処理装置で並列に計算させる。小領域間での中性子束の流出入を解析する際には、前記通信手段により小領域境界の中性子束のデータを演算処理装置間でやり取りすることで、空間的に連続した解析が行われる。

【0005】

【発明が解決しようとする課題】複数の演算処理装置を用いて並列に計算を行わせる際に、演算処理装置間の通信を行うことなく全く独立に計算を進めることができる例はまれであり、通常は演算処理装置間の通信を行いつながら計算を進める。たとえば、4行4列の行列A、Bの掛け算を4台の演算処理装置で実施して4行4列の行列Cを求める場合を考える。A、B、Cの要素をそれぞれ

$A_{11}, B_{11}, C_{11}$  以下のように表記する。

【0006】

【数1】

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \\ b_{31} & b_{32} & b_{33} & b_{34} \\ b_{41} & b_{42} & b_{43} & b_{44} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & c_{13} & c_{14} \\ c_{21} & c_{22} & c_{23} & c_{24} \\ c_{31} & c_{32} & c_{33} & c_{34} \\ c_{41} & c_{42} & c_{43} & c_{44} \end{bmatrix}$$

このとき、4台の演算処理装置のうちの1台においては  
例えば、

【0007】

【数2】

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \\ b_{41} & b_{42} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

のように計算が行われる。

【0008】この例から明らかなように、演算に使う例(1)または(2)については行あるいは列全体にわたる要素のデータが必要である。また、演算の結果として得られるC<sub>11</sub>の方は、各々の演算処理装置に於いては部分的にしかデータが得られない。このことは、例えば次のステップで行列Cと行列Aの掛け算を行う必要が生じたとき、計算で得られた要素だけではデータ不足が生じることを意味する。したがって、A×B=Cの計算を実施した後で残りの部分、上の式で言えば行列Cの少なくとも第1行と第2行のデータ及び第1列と第2列のデータは蓄えられた状態しておかねばならない。

【0009】これらの問題を一般化すると次のようになる。(N×K)個からなる配列X(Nk)があり、これがN台の演算処理装置に分割され、例えば識別番号1の演算処理装置ではX(1)、X(2)、…X(N)、識別番号2の演算処理装置ではX(N+1)、X(N+2)、…X(2N)の計算結果を持っているものとする。この状態からN台の演算処理装置の間で通信を行うことにより、N台の演算処理装置が配列X(Nk)の計算結果を持っている状態を作る操作が必要となることがある。

【0010】このときの通信は1台であることが通信手段上の条件である。すなわち、例えば演算処理装置1から演算処理装置2にデータを転送する際には、演算処理装置2は演算処理装置1からデータを受けとる態勢になければならないのであって、このとき演算処理装置2が他の処理、例えば演算処理装置3にデータを転送しようとしたら演算処理装置4からデータを受けようとしたりすると、通信は失敗して計算は中断することとなる。通信が滞りなく行われるには通信側と受信側の処理がどのように通信の順序を予め決めておく必要がある。

【0011】4台の演算処理装置を使う場合を例にとれば、容易に考えられる方法として次のものが挙げられ

る。以下、表記を簡略化するための演算処理装置1、2、 $\dots$ 、 $N-3$ 、4をそれぞれ#1、#2、#3、#4と書く。

(1) 送信一受信を1つずつ順次行う方法

- |                  |                  |
|------------------|------------------|
| [1] #1の計算結果→#2、  | [2] #1の計算結果→#3、  |
| [3] #1の計算結果→#4、  | [4] #2の計算結果→#1、  |
| [5] #2の計算結果→#3、  | [6] #2の計算結果→#4、  |
| [7] #3の計算結果→#1、  | [8] #3の計算結果→#2、  |
| [9] #3の計算結果→#4、  | [10] #4の計算結果→#1、 |
| [11] #4の計算結果→#2、 | [12] #4の計算結果→#3、 |

を順次実行する。

{0012}にて、[1]、[2]、[3]、 $\dots$ は処理のステップの番号を示す。演算処理装置をN台、1台に割り当てられたデータ量をwとすれば、送信回数

$$2 \times C_t = N(N-1)$$

であり、データ移動量は

$$2w \times C_t = wN(N-1)$$

・データ集約

- |                         |                 |
|-------------------------|-----------------|
| [1] #2の計算結果→#1、         | [2] #3の計算結果→#1、 |
| [3] #4の計算結果→#1、 $\dots$ | を順次実行、          |
| #1に全データが溜る。             |                 |

・全データ配布

- |            |            |                    |
|------------|------------|--------------------|
| [1] #1→#2、 | [2] #1→#3、 | [3] #1→#4、 $\dots$ |
|------------|------------|--------------------|

を順次実行。配列全体を#2、#3、#4に送信する。

{0014}この場合の送信回数は  $2(N-1)$  回、データ移動量は集約時に  $(N-1)w$ 、配布時に  $N(N-1)w$  である。この方法は(1)の方法に比べて送信回数は少ないが、全データ配布時に送信されるデータ

- |                 |            |         |
|-----------------|------------|---------|
| [1] #1の計算結果→#2、 | #3の計算結果→#4 | を同時に実行、 |
| [2] #1の計算結果→#3、 | #2の計算結果→#4 | を同時に実行、 |
| [3] #1の計算結果→#4、 | #2の計算結果→#3 | を同時に実行、 |
| [4] #2の計算結果→#1、 | #4の計算結果→#3 | を同時に実行、 |
| [5] #3の計算結果→#1、 | #4の計算結果→#2 | を同時に実行、 |
| [6] #4の計算結果→#1、 | #3の計算結果→#2 | を同時に実行、 |

{0015}この送信方法によれば、通信が重複することなく、全データが各々の演算処理装置に行き渡る。演算処理装置がN台であれば送信回数は  $2(N-1)$ 、データ移動量は  $2(N-1)w$  である。 $N=4$ であれば送信回数は(2)の方法の1/3、データ移動量は(2)の方法の1/2である。 $N$ が大きくなるとともに通信量は広がる。

{0017}(4) 演算処理装置の Binary treeにより代表の演算処理装置にデータを集めた後、各演算処理装置に配布する。これは(2)の方法を改良したもので、例えば次のように行う。

・データ集約

- |                     |            |         |
|---------------------|------------|---------|
| [1] #2の計算結果→#1、     | #4の計算結果→#3 | を同時に実行、 |
| [2] #3に集約された計算結果→#1 |            |         |
| ・全データ配布             |            |         |
| [3] #1→#3           |            |         |
| [4] #1→#2、          | #3→#4      | を同時に実行、 |

※である。 $N=4$ ならば送信回数は上述の1/2である。この方法によれば、時間はかかるが通信上の混乱は避けられる。なお、 $C_t$ は1個の要素からq個の要素を添え組合せの数を示す。

{0013}(2) 代表の演算処理装置にデータを集めた後、各演算処理装置に配布する。

★タ数が多い点が短所である。

{0015}また、通信の効率化を図った手法として次のものがある。

(3) 演算処理装置の1組1の組み合わせに対して並列・網羅的に通信を行う。これは(1)の方法を改良したもので、例えば次のように行う。

{0018}この方法によれば、演算処理装置がN台であれば、送信回数は  $2 \times \log_2 N$  回、データ送信量は、集約時に  $(N-1)w$ 、配布時に  $Nw \log_2 N$  である。 $N=4$ であれば送信回数は(2)の方法の2/3、データ移動量は(2)の方法の11/15である。 $N$ が大きくなるとともに通信量は広がる。

{0019}(3)の方法は(4)の方法に比べてデータ移動量は少ないが送信回数が多いため、格う配列が小さい場合には適していない。(4)の方法は送信回数は少ないが、データ移動量が多いため、巨大な配列を扱う場合には適していない。

{0020}よって、データ移動量と送信回数がともに最適化された、あらゆる条件に対して適用可能な一般化された手法が必要である。本発明は、このような点を考慮してなされたもので、通信によるデータの検出と並列に行えるようにすることで、演算処理装置間の通信回数およびデータの検出の待ち時間を最小限に抑えて高

連化を図ることができ、並列計算機システム及びその演算処理装置間の通信方法を提供することを目的とする。

【0021】

問題を解決するための手段として、上記目的を達成するため、本発明の請求項1記載の発明は、固有の識別子を有する少なくとも2台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、2<sup>nd</sup> 番の小配列に分割して2<sup>nd</sup> 台の演算処理装置に分配される各演算処理装置で演算処理されたデータ配列を再び1つの配列に集結する際に、2<sup>nd</sup> 台の演算処理装置に識別番号0、1、…、2<sup>nd</sup>-1を付与し、識別番号Nの演算処理装置に対し2番法で表した識別番号Nの2<sup>nd</sup> の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に授受する操作iをi=0からi=m-1まで順次行い、i>nなるjに対しては、操作jの際に、識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作j-1までで得られた演算処理結果を授受することにより2<sup>nd</sup> 台の演算処理装置間でn回の操作でデータ配列を集結させることを特徴とする。

【0022】また、請求項2記載の発明は、固有の識別子を有する(2<sup>nd</sup>+k)台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、(2<sup>nd</sup>+k) 個の小配列に分割して(2<sup>nd</sup>+k) 台の演算処理装置に分配・演算処理されたデータ配列を再び1つの配列に集結する際に、前記(2<sup>nd</sup>+k) 台の演算処理装置に個別記憶手段及び通信手段を備えた(2<sup>nd</sup>-k) 台の演算処理装置を加えた2<sup>nd</sup> 台からなる演算処理装置群を形成し、この演算処理装置群を構成する2<sup>nd</sup> 台の演算処理装置に識別番号0、1、…、2<sup>nd</sup>-1を付与し、識別番号Nの演算処理装置に対し2番法で表した識別番号Nの2<sup>nd</sup> の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に授受する操作iをi=0からi=mまで順次行い、i>nなるjに対しては、操作jの際に、N ≤ 2<sup>nd</sup>+kなる識別番号Nの演算処理装置からはその演算処理装置の演算処理結果及び操作j-1までで得られた演算処理結果を授受し、N > 2<sup>nd</sup>+kなる識別番号Nの演算処理装置からは操作j-1までで得られた演算処理結果を授受することにより(2<sup>nd</sup>+k) 台の演算処理装置において(m+1) 回の操作でデータ配列を集結させることを特徴とする。

【0023】また、請求項3記載の発明は、固有の識別子を有する(2<sup>nd</sup>+k) 台の演算処理装置と、これら各

演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、(2<sup>nd</sup>+k) 個の小配列に分割して(2<sup>nd</sup>+k) 台の演算処理装置に分配・演算処理されたデータ配列を再び1つの配列に集結する際に、この(2<sup>nd</sup>+k) 個のデータ配列に(2<sup>nd</sup>-k) 個の空の小配列を追加することにより前記データ配列を小配列2<sup>nd</sup> 個分の配列に拡張し、前記(2<sup>nd</sup>+k) 台の演算処理装置に個別記憶手段及び通信手段を備えた(2<sup>nd</sup>-k) 台の演算処理装置を加えた2<sup>nd</sup> 台からなる演算処理装置群を形成し、この演算処理装置群を構成する2<sup>nd</sup> 台の演算処理装置に識別番号0、1、…、2<sup>nd</sup>-1を付与し、識別番号Nの演算処理装置に対し2番法で表した識別番号Nの2<sup>nd</sup> の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に授受する操作iをi=0からi=mまで順次行い、i>nなるjに対しては、操作jの際に、識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作j-1までで得られた演算処理結果を授受することにより(2<sup>nd</sup>+k) 台の演算処理装置において(m+1) 回の操作でデータ配列を集結させることを特徴とする。

【0024】また、請求項4記載の発明は、n>mなるn、mについて、固有の識別子を有する(2<sup>nd</sup>+2<sup>nd</sup>) 台の演算処理装置と、これら各演算処理装置に各々対応する個別記憶装置および通信手段とを備え、この通信手段により各演算処理装置間でデータの授受を行う並列計算機システムにおいて、(2<sup>nd</sup>+2<sup>nd</sup>) 個の小配列に分割して(2<sup>nd</sup>+2<sup>nd</sup>) 台の演算処理装置に分配・演算処理されたデータ配列を再び1つの配列に集結する際に、前記(2<sup>nd</sup>+2<sup>nd</sup>) 台の演算処理装置を2<sup>nd</sup> 台からなるグループG<sub>1</sub>と、2<sup>nd</sup> 台からなるグループG<sub>2</sub>に分割し、また前記データ配列を初めの2<sup>nd</sup> 個の小配列からなる配列A、とその後の2<sup>nd</sup> 個の小配列からなる配列A'、の2<sup>nd</sup> に分割し、この配列A<sub>1</sub>、A<sub>2</sub>をそれぞれグループG<sub>1</sub>、G<sub>2</sub>と対応づけて分配、演算処理を行い、グループG<sub>1</sub>の2<sup>nd</sup> 台の演算処理装置に識別番号0、1、…、2<sup>nd</sup>-1を付与し、識別番号Nの演算処理装置に対し2番法で表した識別番号Nの2<sup>nd</sup> の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に授受する操作iをi=0からi=n-1まで順次行い、i>nなるjに対しては、操作jの際に識別番号N、N'の演算処理装置間で各演算処理装置による演算処理結果に加えて操作j-1までで得られた演算処理結果を授受することによりグループG<sub>1</sub>内でデータ配列を集結させる第1の工程と、グループG<sub>2</sub>の2<sup>nd</sup> 台の演算処理装置に換

11

別番号0, 1, ...,  $2^k - 1$ を付与し、識別番号Nの演算処理装置に対し2進法で表した識別番号Nの2<sup>i</sup>の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受信する操作: (i = 0 から i = 0...1)まで繰り返し、i > 0 となる」に対して、操作: j の際に識別番号N, N'の演算処理装置間で各演算処理装置による演算処理結果に追加で操作: (j - 1)まで得られた演算処理結果を送受信することによりグループG<sub>i</sub>内でデータ配列を集結させる第2の工程と、グループG<sub>i</sub>からグループG<sub>i+1</sub>の各演算処理装置に配列A<sub>i</sub>を、グループG<sub>i+1</sub>からグループG<sub>i+2</sub>の各演算処理装置に配列A<sub>i+1</sub>を送信する第3の工程とを有し、第1の工程と第2の工程を並列に実行した後に第3の工程を行なうことにより(2<sup>k</sup> + 2<sup>k-1</sup>) 台の演算処理装置においてデータ配列を集結させることを特徴とする。

【0025】また、請求項5記載の発明は、固有の識別子を有する複数の演算処理装置と、これら各演算処理装置にそれぞれ対応する個別記憶装置および通信手段とを備えた並列計算機システムにおいて、

【数8】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

式の右辺項(但し、 $n_1 > n_2 > n_3 > \dots > n_k$ ,  $n_i \geq 0$ )に分割して

【数9】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

台の演算処理装置に分配、演算処理されたデータ配列を再び1つの記憶に集結する際に、これらの演算処理装置のうち

【数10】

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

台をそれぞれグループG<sub>1</sub>, G<sub>2</sub>, ..., G<sub>k</sub>としてk個のグループに分割するとともに、前記データ配列のうち

【0029】

【数11】

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

台のデータ配列をそれぞれ配列A<sub>1</sub>, A<sub>2</sub>, ..., A<sub>k</sub>としてk個の配列に分割し、このk個の配列とk個のグループG<sub>1</sub>, G<sub>2</sub>, ..., G<sub>k</sub>とを1対1に対応づけて分配、演算処理を行い、1 ≤ p ≤ k なる各pに対して、グループG<sub>p</sub>の(2<sup>n<sub>p</sub></sup>台)の演算処理装置に識別番号0, 1, ...を付与し、識別番号Nの演算処理装置に対し2進法で表した識別番号Nの2<sup>i</sup>の位の数を反転させた番号N'を識別番号とする演算処理装置を対応させ、前記データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置の間で相互に送受信する操

12

作: i を i = 0 から i = 0...1 まで順次行い、i > 0 となる」に対して、操作: j の際に識別番号N, N'の演算処理装置間で各演算処理装置による演算処理結果に追加で操作: (j - 1)まで得られた演算処理結果を送受信することによりグループG<sub>i</sub>内の演算処理装置でデータ配列A<sub>i</sub>を集結させるグループ内工程hを実行し、グループ内工程(h - 1)が終了した後、グループG<sub>i</sub>の演算処理装置から配列A<sub>i</sub>の演算結果をグループG<sub>i+1</sub>の演算処理装置に送信するグループ間工程hを実行し、次にグループG<sub>i</sub>の各演算処理装置に集結された配列A<sub>i</sub>の演算結果を、グループG<sub>i</sub>の演算処理装置からq > p なる全てのqに対しグループG<sub>i</sub>に属する各演算処理装置に送信するとともに、グループG<sub>i</sub>の演算処理装置から、グループG<sub>i</sub>自身の演算結果である配列A<sub>i</sub>及びグループG<sub>i+1</sub>の演算処理装置から受信した配列A<sub>i+1</sub>, ..., A<sub>k</sub>の演算結果をグループG<sub>i+1</sub>の演算処理装置に送信するグループ間工程pを、p = k - 1 から p = 2 までpに關して順次実行することにより、

【0027】

【数12】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

台の演算処理装置においてデータ配列を集結させることを特徴とする。

【0028】なお、この際には、k個のグループ内工程1, 2, ..., kを並列に実行し、1 ≤ s ≤ k - 1 なるsに対して、グループ内工程sが終了した時点で順次グループ間工程(s + 1)を実行することで、全体の通信に要する時間をさらに短縮することができる。

30

【0029】また、請求項6記載の発明は、請求項4または5記載の並列計算機システムを用いて演算処理装置のグループ間でデータ交換を行う場合、2<sup>k</sup>台の演算処理装置からなるグループG<sub>i</sub>で集結され共有されているデータ配列A<sub>i</sub>と、2<sup>k</sup>台の演算処理装置(p > i)からなるグループG<sub>j</sub>で集結され共有されているデータ配列B<sub>j</sub>とを、グループG<sub>i</sub>, G<sub>j</sub>間で相互に送受信する際に、グループG<sub>i</sub>のなかから選択される2<sup>i</sup>台の演算処理装置をグループG<sub>i</sub>の各演算処理装置と1対1に対応させてグループG<sub>j</sub>の各演算処理装置にデータ配列A<sub>i</sub>を送信する操作を並列に実施するとともに、グループG<sub>j</sub>を、それぞれが2<sup>j-i</sup>台の演算処理装置からなる小グループα<sub>1</sub>, α<sub>2</sub>, ..., α<sub>r</sub> (r = 2<sup>j-i</sup>)に分割して、各々の小グループとグループG<sub>i</sub>のr個の演算処理装置とを1対1に対応させ、小グループα<sub>q</sub>のなかから選択される1台の演算処理装置に対して、小グループα<sub>q</sub>に対応するグループG<sub>i</sub>の演算処理装置からデータ配列B<sub>j</sub>を送信した後、小グループα<sub>q</sub>の演算処理装置間でデータ配列B<sub>j</sub>を送受信する操作: i を、1 ≤ i ≤ r となるiに關して並列に実行するすることにより、2<sup>k</sup>台の演算処理

50

装置と2' 台の演算処理装置にデータ配列Aとデータ配列Bを共有させることを特徴とする。

【0030】また、請求項1記載の発明は、請求項1乃至前記のいずれか記載の並列計算機システムを用いて2台の演算処理装置間でデータを交換する工程は、演算処理装置の識別番号の大きい方から小さい方にデータを送る第1の送信工程と、演算処理装置の識別番号の小さい方から大きい方にデータを送る第2の受信工程とからなり、この第1の送信工程と第2の受信工程のうちから選択される1工程を先に行った後、続いて他の1工程を行うことを特徴とする。

【0031】

【発明の実施の形態】本発明の実施の形態について、以下、図面を参照して説明する。図1は並列計算機システムの構成例を示すブロック図である。ここに示した並列計算機システムは、1台のホストの計算機1と8台の演算処理装置2-1、2-2、2-3、2-4、2-5、2-6、2-7、2-8で構成されている。ホストの計算機には記憶装置3と通信手段4、演算処理装置2-1、2-2、2-3、2-4、2-5、2-6、2-7、2-8の各々には、制御記憶装置5-1、5-2、5-3、5-4、5-5、5-6、5-7、5-8と通信手段6-1、6-2、6-3、6-4、6-5、6-6、6-7、6-8が備えられている。例えば、ホストの計算機で読み込んだ入力データ等は、通信手段4から通信手段6-1、6-2、6-3、6-4、6-5、6-6、6-7、6-8を通じて演算処理装置に送られる。演算処理装置2-1、2-2、2-3、2-4、2-5、2-6、2-7、2-8では各々割り当てられた領域の計算を行い、必要に応じて演算処理装置間の通信によりデータの授受を行う。

【0032】図1に示した並列計算機システムの構成に基づき、本発明にかかる並列計算機システムの第1の実施の形態について説明する。図2は本実施の形態における並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。

【0033】演算処理装置2-1、2-2、2-3、2-4、2-5、2-6、2-7、2-8の識別番号をそれぞれ、1、2、3、4、5、6、7、8とし、これらを2進法の3桁の数字として表示するとそれぞれ000、001、010、011、100、101、110、111となる。8×N個のデータからなる配列AがK個のデータからなる8個の小配列A<sub>1</sub>、A<sub>2</sub>、A<sub>3</sub>、A<sub>4</sub>、A<sub>5</sub>、A<sub>6</sub>、A<sub>7</sub>、A<sub>8</sub>に分割されて、8台の演算処理装置2-1、2-2、2-3、2-4、2-5、2-6、2-7、2-8に割り当てられている。それぞれの演算処理装置で割り当てられた小配列のデータに関する演算処理を行った後、配列Aの要素を全ての演算処理装置において処理することを考える。なお、図2において各演算処理装置にかかれた0または1はそれぞれ保持された小配列を示しており、1は計算結果が未入力の状態を、0は計算結果が入力済みの状態を示す。

【0034】第1ステップとして、2'の位の数を反転させた値を識別番号としてもつ演算処理装置との間でデータを交換する。例えば演算処理装置0100は演算処理装置1001、演算処理装置0111は演算処理装置2010とN個のデ

ータを交換する。各演算処理装置に2'個の要素が備えられている。

【0035】第2ステップでは、2'の2位の数を反転させた値を識別番号としてもつ演算処理装置との間でデータを交換する。例えば演算処理装置0100は演算処理装置2010、演算処理装置0111は演算処理装置1001とデータを交換する。この時、例えば演算処理装置010から演算処理装置20への通信では、演算処理装置0自身による演算結果の他に第1ステップで演算処理装置1から受信したデータを含む2'個のデータを通信する。これにより各演算処理装置に4'個の要素が備えられる。

【0036】最後に第3ステップとして、2'の3位の数を反転させた値を識別番号としてもつ演算処理装置との間でデータを交換する。例えば演算処理装置0100は演算処理装置4100、演算処理装置0111は演算処理装置7011と4'個のデータを交換する。各演算処理装置に8'個の要素が備えられ、操作が完了する。

【0037】以上述べた通信方法は演算を2' = 8個に分割した場合であり、この時のステップ数は3である。同様に、演算を2' = 16個に分割した6台の演算処理装置において通信を行なう場合には、上述した8分割の場合に比べてさらに1ステップが必要となり、全部で4ステップとなる。一般に、演算をN個に分割したN台の演算処理装置において通信を行う場合は、上述の方法を採用して、ステップ数  $\log_2 N$  で通信が完了する。

【0038】本実施の形態の作用効果について以下に概説する。例えば配列の大きさをM[word]、演算処理装置の台数をNとし、配列全体がKに分割されて各演算処理装置に置かれているものとする。Kの値としては並列計算で最も一般的な条件である2のべき乗の場合、つまりK = 2' と表される場合について考える。この状態から、演算処理装置間の通信によって演算処理装置全部が配列全体についてデータを把握している状況を作り出すのにかかる時間について考察する。一般にデータを送受するのにかかる時間丁は

$$T = A + B \times W \quad \dots \dots \dots (1)$$

と表せる。ここで、Aは通信準備に要する時間で、送受するデータ量に関わらず1回の通信に必ず必要となる時間である。Aの値はデータ量に依らない。B×Wはデータ量に比例する項であり、Wがデータ量（word数）、Bが1word当たりの転送時間である。

【0039】データの搬送のステップ数は  $\log_2 K = \log_2 N$  である。各ステップで演算処理装置間で通信と受信が1回ずつ行われる。第mステップで搬送されるデータ量は  $(M/K) \times 2^m$  [word] である。データ量M[word]のデータを全演算処理装置において集約せしめるのに必要な送受信の総数は各演算処理装置当たり2'回であり、送受信する総データ量は

$$M \times 2^m \times 2^m$$



【図13】

$$\sum_{i=1}^n \frac{M}{K} 2^i = \frac{M}{K} 2(2^n - 1) = 2M(1 - \frac{1}{K}) \quad [\text{word}]$$

である。よって、本発明を適用した場合の全通信時間 T は、

$$T(K) = 2A \log_2 K + 2M(1 - 1/K)B \quad (3)$$

となる、

【0041】比較のため、従来法、例えば Binary tree の方式で 1 台の演算処理装置に全データを集めておき、同様に Binary tree の方式で全演算処理装置にデータを配信する場合の通信時間を次に示してある。全データ×10

$$T_1(K) = A \log_2 K + M(1 - 1/K)B \quad (4)$$

となる、

【0042】代表演算処理装置から各演算処理装置にデータを配布する際のステップ数は  $\log_2 K$  で、演算処理装置あたり通信回数も最大で  $\log_2 K$  回である。ただ ★

$$T_1(K) = A \log_2 K + M \log_2 K B \quad (5)$$

となる。したがって、全通信時間 T<sub>1</sub> = T<sub>1</sub> + T<sub>2</sub> は

$$T_1(K) = 2A \log_2 K + M(1 - 1/K + \log_2 K)B \quad (5)$$

となる、

【0043】図5及び図6のグラフは、機械に演算処理装置台数、ネットワークに要する時間によって、演算処理装置台数増加に伴う通信時間の増加の関係を示しており、従来の Binary tree の通信方式による (5) 式の関係と、本実施の形態により通信を効率化した (3) 式の関係を比較して示している。このグラフ中の曲線のうち実線で示した符号 10a、10b が本実施の形態の (2) 式の場合、破線で示した符号 11a、11b が従来の (5) 式の場合を示している。図3に示した符号 10a、11a を付した曲線は、通信されるデータ量が少なく、(1) 式の A (通信立ち上げ時間) が全通信時間 T<sub>1</sub> のほぼ半分を占める状況を示す。また図4に示した符号 10b、11b を付した曲線は、通信されるデータ量が多く、(3) 式の A (通信立ち上げ時間) が全通信時間 T<sub>1</sub> に比べて十分小さい状況を示している。このグラフからも明らかなように、本実施の形態によれば、演算処理装置の台数が少数の場合、多数の場合何れも従来の方法より通信に要する時間を少なくすることができ、すなわち、本実施の形態により、データの授受の際の待ち時間を最小限に抑え、計算の高速化を図ることができ、

【0044】なお、例えば演算処理装置の台数が 16 台からなる並列計算機システムにおいて、その内の 8 台の演算処理装置の間で上述の 3 ステップからなる分割の分割、演算結果を行うなど、複数の演算処理装置のうち 2 の装置の台数だけ抜き出してこれらに通信制御用の識別番号を付し、この台数に達して上述した方法で分割の分割分配、演算処理を行なうものとしてもよい。

【0045】上記第 1 の実施の形態においては、関係する演算処理装置の台数が 2 の事象であることを前提としている。一般的な条件として演算処理装置の台数が 2 の事象でない場合、すなわち台数が  $2^k + k$  等として表

すデータを 1 台の演算処理装置に集めるのに要する送受信の回数は、代表演算処理装置において  $n = \log_2 K$  回である。また、第 m ステップ (m ≤ n) で送信されるデータ量は  $(M/K) \times 2^{m-1}$  [word] である。よって、代表演算処理装置に全データを集めるのにかかる時間 T<sub>2</sub> は

★し、各ステップ毎に M [word] のデータが送信される。よって、各演算処理装置にデータを配布する際にかかる時間 T<sub>2</sub> は

される場合にも拡張したのが以下詳述する第 2 の実施の形態である。

【0046】本発明にかかる並列計算機システムの第 2 の実施の形態について説明する。ここでは、例えば並列計算の配列を 8 分割して、8 台の演算処理装置 (識別番号を 0、1、…、7 とする。) に割り当てる場合について説明する。図 8 は本実施の形態における並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。この際のデータ処理には、図 8 台の演算処理装置のほかに 2 台の演算処理装置 (識別番号を 8、7 とする。) を用いることとする。

【0047】第 1 ステップとして、2<sup>6</sup> の位の数を反転させた数を識別番号としてもつ演算処理装置 0 の間でデータを交換する。例えば演算処理装置 0 (000) は演算処理装置 1 (001) と、演算処理装置 3 (011) は演算処理装置 2 (010) と、それぞれ n 個のデータを交換する。演算処理装置 6 (110) と演算処理装置 7 (111) は交換すべきデータがないので休止する。この時点で、演算処理装置 0～6 に 2 n 個のデータが集められる。

【0048】第 2 ステップでは、2<sup>5</sup> の位の数を反転させた数を識別番号としてもつ演算処理装置 0 の間でデータを交換する。例えば演算処理装置 0 (000) は演算処理装置 6 (110) とのデータ交換となるが、この時点で演算処理装置 6 (110) は送信すべきデータがないので、演算処理装置 4 からデータを受信するのみとする。このデータ交換により、演算処理装置 0～3 に 4 n 個のデータが、演算処理装置 4～7 に 2 n 個のデータが集められる。

【0049】第 3 ステップでは、2<sup>4</sup> の位の数を反転させた数を識別番号としてもつ演算処理装置 0 の間でデータを交換する。例えば演算処理装置 0 (000) は演算処理装置 2 (010) との交換である。演算処理装置 8 から演算

処理装置2へは2台のデータ、演算処理装置2から演算処理装置6へは4台のデータを送信する。このようにして6台のデータが6台の演算処理装置全てに伝送される。

【0050】本実施の形態においては、一般に $(2^k + k)$ 台の演算処理装置に対して、 $(2^k - k)$ 台の演算処理装置を加えた $2^{k+1}$ 台の演算処理装置群を構成し、この演算処理装置群に対して上述の第1の実施の形態で述べたステップにより並列計算を行うものとする。これにより、2の乗数ではない台数の演算処理装置に対して2の乗数の場合に準じた構成とすることで、上記第1の実施の形態と同様の作用効果を得ることができる。

【0051】次に本発明にかかる並列計算機システムの第3の実施の形態を説明する。本実施の形態における演算処理装置間の通信方法について、例として、配列を6個の小配列分割して6台の演算処理装置（識別番号0、1、…、5）に割り当てている場合について説明する。まず前記配列を小配列と分割し、拡張した部分には0を埋める。例えば12個の要素からなる配列 $\{3, 1, 4, 1, 5, 9, 6, 5, 3, 5, 8, 0, 0, 0, 0\}$ を拡張からなる配列 $\{0, 0, 0, 0\}$ を追加して、16の要素からなる配列 $\{3, 1, 4, 1, 5, 9, 6, 5, 3, 5, 8, 0, 0, 0, 0\}$ とする。演算処理装置としては前記6台の演算処理装置のほかに2台の演算処理装置（識別番号8、7とする）を加えた8台の演算処理装置を用いる。この場合は、上記第1の実施の形態において述べた手順により、8台の演算処理装置間で通信を行い、データを交換する。

【0052】本実施の形態においては、一般に $(2^k + k)$ 台の演算処理装置に対して、 $(2^k - k)$ 台の演算処理装置を加えた $2^{k+1}$ 台の演算処理装置群を構成し、また配列についてもその要素を $2^{k+1}$ 個に拡張して各演算処理装置に分配し、上記第1の実施の形態と同様の方法で並列計算及びデータの通信を行うものとする。これにより、2の乗数ではない台数の演算処理装置に対して2の乗数の場合に準じた構成とすることで、上記第1の実施の形態と同様の作用効果を得ることができる。

【0053】次に、本発明にかかる並列計算機システムの第4の実施の形態について説明する。第2及び第3の実施の形態は、配列の分割数が2の乗数でない場合、すなわち $(2^k + k)$ 台に分割される場合について、 $2^{k+1}$ 台の演算処理装置によってデータ配列を1個に集約する方法について述べたものである。これに対し本実施の形態は、配列の分割数が $2^k + 2^m$ （ $n \geq m$ ）である場合に対し、 $(2^n - 2^m)$ 台の演算処理装置で処理するものである。

【0054】本実施の形態における並列計算機システムの演算処理装置間の通信方法として、ここではまず例として、配列を8分割して8台の演算処理装置（識別番号0、1、…、7）に割り当てられている場合について説明する。図8はこの場合の演算処理装置間通信方法を時系列

で示すチャートである。

【0055】まず、8台の演算処理装置を2つのグループに分割する。演算処理装置グループ1は識別番号0〜3の4台で構成される。演算処理装置グループ2は識別番号4〜7の4台で構成される。次に、演算処理装置グループ1の4台間、および演算処理装置グループ2の4台間で、上述の第2の実施の形態における手順により、各々のグループでデータを連結させる。図8における第1及び第2ステップがこれに相当する。

【0056】この後、グループ1とグループ2でデータ交換を次の手順で行う。

[1] 演算処理装置4→演算処理装置0、演算処理装置5→演算処理装置2を同時並列に実施。（図6の第3ステップに相当。）

[2] 演算処理装置1→演算処理装置4、演算処理装置3→演算処理装置6を同時並列に実施。（図6の第4ステップに相当。）

[3] 演算処理装置0→演算処理装置1、演算処理装置2→演算処理装置3を同時並列に実施。（図6の第5ステップに相当。）

この方法により、8台の演算処理装置によってデータ配列を集約させることができる。

【0057】また、本実施の形態のもう一つの例として、配列を16分割して16台の演算処理装置（識別番号0、1、…、15）に割り当てられている場合について説明する。図7はこの場合における演算処理装置間通信方法を時系列で示すチャートである。

【0058】まず、16台の演算処理装置を2つのグループに分割する。演算処理装置グループ1は識別番号0、1、…、7の8台で構成される。演算処理装置グループ2は識別番号8、9の8台で構成される。次に演算処理装置グループ1の8台の演算処理装置間、および演算処理装置グループ2の8台の演算処理装置間で、上記第2の実施の形態において述べた方法により、各々のグループでデータを連結させる。図7における第1、第2及び第3ステップがこれに相当する。

【0059】この後、グループ1とグループ2でデータ交換を次の手順で行う。

[1] 演算処理装置8→演算処理装置0、演算処理装置9→演算処理装置4を同時並列に実施。（図7の第4ステップに相当。）

[2] 演算処理装置3→演算処理装置8、演算処理装置7→演算処理装置6、演算処理装置0→演算処理装置2、演算処理装置4→演算処理装置9を同時並列に実施。（図7の第5ステップに相当。）

[3] 演算処理装置0→演算処理装置1、演算処理装置2→演算処理装置3、演算処理装置4→演算処理装置5、演算処理装置6→演算処理装置7を同時並列に実施。（図7の第6ステップに相当。）

【0060】この方法により、8台の演算処理装置によ

ってデータ配列を連結させることができる。なお、グループ2からグループ1に送信されたデータのグループ2内の分配は Binary Tree の方式によっている。

【0001】以下、本発明にかかる並列計算機システムの第1の実施形態について説明する。本実施形態における演算処理装置間の通信方法は、上記第1の実施形態の通信方法を一般化したものである。以下、例として図1を22分けて22台の演算処理装置（識別番号0、1、…、21）に割り当てている場合について説明する。図4及び図9はこの配列に分類した場合における演算処理装置間通信方法を時系列で示すチャートである。図8において第1ステップから第4ステップまでを、図9において第5ステップから第8ステップまでを示した。

【0002】 $22 = 2^4 + 2^2 + 2^0$  であるから、まず、演算処理装置を次の3グループに分ける。

グループ1： 識別番号0、1、…、19の演算処理装置（19台）

グループ2： 識別番号16、17、18、19の演算処理装置（4台）

グループ3： 識別番号20、21の演算処理装置（2台）

【0003】次に、演算処理装置グループ1の10台間、演算処理装置グループ2の4台間、および演算処理装置グループ3の2台間で、上記第1の実施形態の方法により各々のグループでデータを連結させる。これは図8に示した第1ステップから第4ステップまでが相当する。

【0004】この後は、上記第2ないし第3の実施形態において説明した方法と同様の手順により、データの

グループ1の小グループ1： 演算処理装置0、1、2、3

小グループ2： 演算処理装置4、5、6、7

小グループ3： 演算処理装置8、9、10、11

小グループ4： 演算処理装置12、13、14、15

とする。

【0005】この各小グループから1台ずつ演算処理装置を選択する。ここでは演算処理装置0、4、8、12を選択する。この4台の演算処理装置に対して、それぞれグループ2の演算処理装置16、17、18、19から、グループ2及びグループ3に関して連結されたデータを送信する。これは図9に示した第5ステップに相当する。

【0006】次に、グループ1の各小グループにおいて、従来の Binary Tree の方式で演算処理装置間でグループ2、3に関するデータの送受信を行ない、小グループの演算処理装置においてグループ1、2、3のデータを連結させる。例えば小グループ1においては演算処理装置0から演算処理装置2に対してデータを送信し、次に演算処理装置0、2からそれぞれ演算処理装置1、3に対してデータの送信を行う。他の小グループにおいても同様である。これは図9に示した第7ステップ及び第8ステップに相当する。こうして、全ての22台の演算処理装置において22個のデータ配列の連結を完了する。

\* グループ間交換を行う。以下そのデータの通信方法を順を追って説明する。まず、第2ステップでグループ2においてデータの連結が終了するが、その時点で既にグループ3のデータの連結は完了しているから、次のステップとして、グループ2の演算処理装置16、18とグループ3の演算処理装置19、20との間でそれぞれデータの交換が行なわれる。これは図8に示したグループ2とグループ3における第3ステップに相当する。この時点で、グループ3の全ての演算処理装置にはグループ2及びグループ3におけるデータがすべて格納された状態となる。

【0005】次に、グループ2及びグループ3の全てのデータが格納されたグループ2の演算処理装置16、18から、それぞれグループ2の演算処理装置17、19に対してグループ3より受信したデータが送信される。これは図8に示したグループ2における第4ステップに相当する。

【0006】グループ1においては第4ステップで各演算処理装置間でデータの連結が終了するが、次のステップとして、グループ1と、グループ2、3との間でデータの送受信を行う。まず、グループ1の演算処理装置0、1、2、3、4、5から、それぞれグループ2、3の演算処理装置16、17、18、19、20、21に対してデータが送信される。これによりグループ2、3においてはグループ1、2、3の22台の全ての演算処理装置のデータの連結が完了する。これは図9に示した第8ステップに相当する。

【0007】次に、グループ1の10台の演算処理装置を4つの小グループに分割する。すなわち、

グループ1の小グループ1： 演算処理装置0、1、2、3

小グループ2： 演算処理装置4、5、6、7

小グループ3： 演算処理装置8、9、10、11

小グループ4： 演算処理装置12、13、14、15

【0008】一様に、2の乗数では表されない任意の演算処理装置におけるデータ配列は、以上説明した方法によって連結させることができる。まず、各個の整数  $n_1, n_2, n_3, \dots, n_k$ 、 $(n_1, n_2, n_3, \dots, n_k)$  を用いて、並列計算機システムの演算処理装置の台数を

【0009】

【数1】

$$2^{n_1} + 2^{n_2} + 2^{n_3} + \dots + 2^{n_k}$$

と表す。また、データ配列をこの台数と同数の小配列に分割し、各演算処理装置に分割して演算処理を行なうものとする。並列計算機システムの演算処理装置のうち、

【0010】

【数1】

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

台をそれぞれグループG<sub>1</sub>、G<sub>2</sub>、…、G<sub>k</sub>として、並列計算機システムの演算処理装置をk個のグループに分

測する。同時にデータ配列の小配列の

{007B}

{数1B}

$$2^{n_1}, 2^{n_2}, 2^{n_3}, \dots, 2^{n_k}$$

個をそれぞれ配列A<sub>1</sub>, A<sub>2</sub>, ..., A<sub>k</sub>としてk個の配列に分割する。

{0074} 次に、1 ≤ p ≤ k なるすべてのpに対して、以下の「」内に定義する操作（以下、グループ内工程pという。）を行う。但し、グループ内工程1, ..., kは並列して行うこととする。

{0075} i グループG<sub>i</sub>の（2のn<sub>i</sub>乗）個の演算処理装置に識別番号0, 1, ..., (2のn<sub>i</sub>乗-1)を付与する。次に、0 ≤ q ≤ p-1なるqに対し、以下の「」内に定義する操作qを、q=0からq=p-1まで順次行う。

① 識別番号Nの演算処理装置に対し、2進法で表した識別番号Nの2<sup>i</sup>の位を反転させた番号N'を識別番号とする演算処理装置を対応させ、データ配列の演算処理結果を識別番号Nの演算処理装置と識別番号N'の演算処理装置との間で相互に送受信する。但し、q>0なるqに対しては、操作qの際に、識別番号N、N'の演算処理装置間で、各演算処理装置による演算処理結果に加えて操作（q-1）までで得られた演算処理結果を含むて送受信することとする。

② この操作により、グループG<sub>i</sub>の（2のn<sub>i</sub>乗）台の演算処理装置で、データ配列の集結を行う。

グループの設定方法により、グループ内工程1, ..., kを並列に行なったとき、グループ内工程kが最初に終了し、以下、グループ内工程（k-1）, ..., 2, 1の順に終了する。このことを考慮して、以下の「」に定義する操作（以下、グループ内工程pという。）を、p=k-1からp=1までpに関して順次に行うこととする。

{0076} i { グループ内工程pが終了した後、グループG<sub>i</sub>の各演算処理装置に集結された配列A<sub>i</sub>のデータを、グループG<sub>i</sub>の演算処理装置から、グループG<sub>1</sub>, ..., G<sub>i-1</sub>, G<sub>i+1</sub>, ..., G<sub>k</sub>に属する全ての演算処理装置に送信する。すなわち、グループG<sub>i</sub>に属する（2のn<sub>i</sub>乗）台の演算処理装置のうち

{0077}

{数17}

$$2^{n_{p+1}} + \dots + 2^{n_k}$$

台を選択して、これら選択された演算処理装置とグループG<sub>1</sub>, ..., G<sub>i-1</sub>, G<sub>i+1</sub>, ..., G<sub>k</sub>に属する演算処理装置とを1対1に対応させ、グループG<sub>i</sub>からグループG<sub>1</sub>, ..., G<sub>i-1</sub>, G<sub>i+1</sub>, ..., G<sub>k</sub>への配列A<sub>i</sub>のデータ送信を行う。次に、グループG<sub>1</sub>, ..., G<sub>i-1</sub>からグループG<sub>i</sub>へのデータの送信を行う。（2のn<sub>i</sub>乗）台の演算処理装置からなるグループG<sub>i</sub>を、そ

れぞれが

{0078}

{数18}

$$2^{n_{p+1}} + \dots + 2^{n_k}$$

台の演算処理装置からなる小グループα<sub>1</sub>, ..., α<sub>i</sub>に分割する。この小グループの数tは、

{0079}

{数19}

$$r = 2^{n_{p+1}}$$

である。ここで、グループG<sub>i+1</sub>に属する演算処理装置をb<sub>1</sub>, ..., b<sub>t</sub>と表記する。グループG<sub>i</sub>の小グループα<sub>1</sub>, ..., α<sub>t</sub>と、グループG<sub>i+1</sub>に属する演算処理装置をb<sub>1</sub>, ..., b<sub>t</sub>とを1対1に対応させて、グループG<sub>i+1</sub>の演算処理装置b<sub>j</sub>から対応する小グループα<sub>j</sub>のうちから選択されたt台の演算処理装置a<sub>j</sub>に、グループG<sub>i+1</sub>において集結された配列A<sub>i+1</sub>のデータを送信する操作を、1 ≤ j ≤ tなる全てのjについて並列に行う。このとき、p=k-1の場合、演算処理装置b<sub>j</sub>からa<sub>j</sub>へは、グループG<sub>i+1</sub>, ..., G<sub>k</sub>より受信したデータ配列A<sub>i+1</sub>, ..., A<sub>k</sub>を含めて送信するものとする。

{0080} この後、各小グループα<sub>j</sub>において、演算処理装置a<sub>j</sub>からa<sub>j</sub>以外の全ての演算処理装置に対して、従来のBinary Tree方式でデータの送信を行う。これにより、グループG<sub>p</sub>の全ての演算処理装置に対してデータ配列A<sub>1</sub>, ..., A<sub>k</sub>に関するデータ配列の集結が完了する。}

この方法により、一般に複数台の演算処理装置によって各演算処理装置において分割され並列計算されたデータ配列を、効率よく集結させることができるから、計算の高速化を図ることができる。

{0081}

{発明の効果} 以上説明したように本発明によれば、並列計算システム中の演算処理装置間の通信方法の効率をより向上させることにより、データの授受の際の待ち時間を最小限に抑えることができるから、並列計算システムにおいて実施される大規模な計算の高速化を図ることができる。

{図面の簡単な説明}

{図1} 本発明の第1の実施形態における並列計算システムの構成を示すブロック図である。

{図2} 本発明の第1の実施形態にかかる並列計算システム中の演算処理装置間の通信方法を時系列で示すチャートである。

{図3} 送信されるデータ量が少ない場合の本発明の第1の実施形態及び従来の通信方法を用いた場合の演算処理装置と通信時間の相関を示すグラフである。

{図4} 送信されるデータ量が多い場合の本発明の第1

の実施形態及び従来の通信方法を用いた場合の演算処理回数と通信時間の相関を示すグラフである。

【図5】本発明の第2の実施形態にかかる並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。

【図6】本発明の第2の実施形態にかかる並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。

【図7】本発明の第4の実施形態にかかる並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。

【図8】本発明の第5の実施形態にかかる並列計算機

システムの演算処理装置間の通信方法を時系列で示すチャートである。

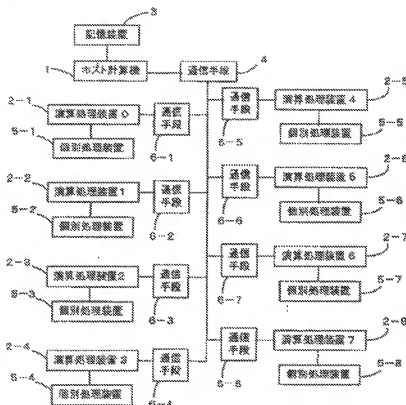
【図9】本発明の第5の実施形態にかかる並列計算機システムの演算処理装置間の通信方法を時系列で示すチャートである。

【符号の説明】

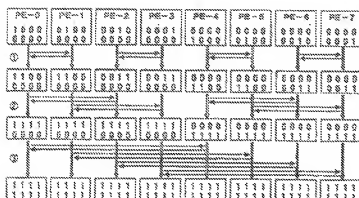
1…ホスト計算機、2…1…演算処理装置、3…記憶装置、4…通信手段、5-1…個別処理装置、5-1…通信手段

10a、10b…本発明の第1の実施形態における演算処理装置の台数と通信に要する時間の関係を示す曲線

【図1】

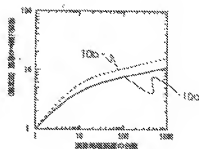


〔図2〕

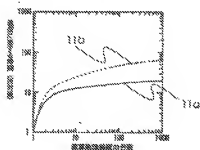


$a_0, a_1, a_2$  : データ入力の状態  
 $b_0, b_1, b_2$  : データ出力の状態

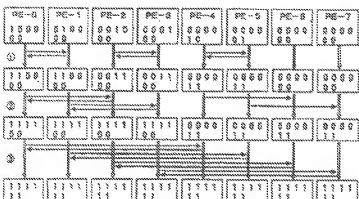
〔図3〕



〔図4〕

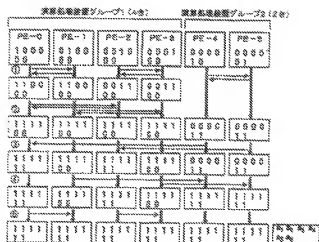


〔図5〕

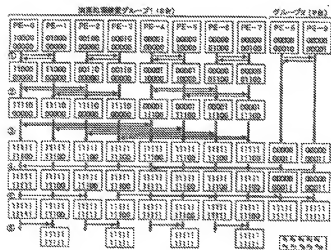


$a_0, a_1, a_2$  : データ入力の状態  
 $b_0, b_1, b_2$  : データ出力の状態

【図3】



【図7】







【図 9】

